

Ejercicio 2: Errores de Medición y Supuestos CEV

Errores en Variables Clásicos (Classical Errors-in-Variables)

Contexto del Problema

Queremos estimar la siguiente ecuación:

$$\text{tvhours}^* = \beta_0 + \beta_1 \text{age} + \beta_2 \text{age}^2 + \beta_3 \text{motheduc} + \beta_4 \text{fatheduc} + \beta_5 \text{sibs} + u \quad (1)$$

donde:

- tvhours^* = horas **verdaderas** de ver televisión por semana (variable **no observable**)
- tvhours = horas **reportadas** en la encuesta (variable **observable**)
- u = error del modelo (término de perturbación estocástica)

Problema fundamental: No observamos tvhours^* , solo observamos tvhours .

Modelo de error de medición

Asumimos que la relación entre el valor verdadero y el observado es:

$$\text{tvhours} = \text{tvhours}^* + e \quad (2)$$

donde e es el **error de medición**.

Pregunta a: ¿Se cumplen los supuestos CEV?

Respuesta

Es **poco probable** que se cumplan perfectamente los supuestos de errores en variables clásicos (CEV) en esta aplicación. Analicemos cada supuesto:

1. Media cero: $E[e] = 0$

Supuesto: En promedio, no hay sesgo sistemático en el reporte.

¿Se cumple? PROBABLEMENTE NO

Razones:

- **Sesgo de deseabilidad social:** Ver mucha televisión puede percibirse negativamente. Los padres o niños pueden subreportar sistemáticamente las horas de TV para dar una mejor impresión.
- **Sesgo de memoria selectiva:** Las personas tienden a recordar mejor ciertos eventos que otros. Pueden olvidar horas "pasivas" de TV mientras hacen otras actividades.
- **Redondeo sistemático:** Si todos redondean hacia abajo (ej. ünas 2 horas cuando en realidad son 2.7), entonces $E[e] < 0$.

Implicación formal: Si $E[e] \neq 0$, entonces:

$$E[\text{tvhours}] = E[\text{tvhours}^*] + E[e] \neq E[\text{tvhours}^*] \quad (3)$$

El nivel promedio reportado difiere sistemáticamente del verdadero, pero esto **NO sesga los coeficientes de pendiente** si los demás supuestos se mantienen.

2. Independencia del valor verdadero: $\text{Cov}(\text{tvhours}^*, e) = 0$

Supuesto: El error de medición no se correlaciona con el valor verdadero.

¿Se cumple? DUDOSO

Razones por las que podría no cumplirse:

1. Heterogeneidad en la precisión del reporte:

- Niños que ven mucha TV (valores altos de tvhours^*) pueden perder la cuenta más fácilmente
- El error podría aumentar con el nivel de TV: $\text{Var}(e|\text{tvhours}^*) = \sigma_e^2(\text{tvhours}^*)$

2. Sesgo diferencial de报告:

- Quien ve mucha TV puede sentir más presión para subreportar (sesgo social más fuerte)
- Esto crearía: $E[e|\text{tvhours}^* = \text{alto}] < 0$ y $E[e|\text{tvhours}^* = \text{bajo}] \approx 0$
- Por lo tanto: $\text{Cov}(\text{tvhours}^*, e) < 0$

Implicación: Si este supuesto falla, el error de medición deja de ser clásico y puede generar sesgos más complejos en los estimadores.

3. Independencia de las X's: $\text{Cov}(X_j, e) = 0$ para toda X_j

Supuesto: El error de medición no se correlaciona con las variables explicativas.

¿Se cumple? PROBABLEMENTE NO

Razones específicas en este contexto:

1. Correlación con educación parental (motheduc, fatheduc):

- Padres con mayor educación pueden:
 - Ser más conscientes del estigma de ver mucha TV
 - Subreportar más las horas de TV de sus hijos
 - Tener mejor memoria/registros de actividades
- Esto implica: $\text{Cov}(\text{motheduc}, e) < 0$

2. Correlación con edad (age):

- Niños más pequeños: reportes hechos por padres (posiblemente más precisos)
- Niños mayores: auto-reportes (posiblemente menos preciso o con más sesgo social)
- Adolescentes pueden ser menos cooperativos en encuestas
- Esto puede crear: $\text{Cov}(\text{age}, e) \neq 0$

3. Correlación con número de hermanos (sibs):

- Familias grandes: más caótico, más difícil hacer seguimiento
- Mayor error de medición en familias numerosas
- Posible: $\text{Cov}(\text{sibs}, e) > 0$ en términos de varianza del error

Implicación crítica: Si $\text{Cov}(X_j, e) \neq 0$, entonces el error de medición se comporta como una variable omitida en la ecuación transformada, lo que puede sesgar los estimadores de todos los coeficientes.

Conclusión sobre el cumplimiento de CEV

Evaluación General

Los supuestos CEV **probablemente NO se cumplen perfectamente** en esta aplicación debido a:

1. Sesgo social de deseabilidad $\rightarrow E[e] \neq 0$
2. Error diferencial según nivel de TV $\rightarrow \text{Cov}(\text{tvhours}^*, e) \neq 0$
3. Correlación con educación parental y edad $\rightarrow \text{Cov}(X_j, e) \neq 0$

Sin embargo, los supuestos CEV pueden ser una **aproximación razonable** si:

- Las violaciones son pequeñas en magnitud
- Los sesgos se cancelan parcialmente
- Las correlaciones son débiles

Pregunta b: ¿Qué exigen los supuestos CEV en esta aplicación?

Definición formal de los supuestos CEV

Los supuestos de **errores en variables clásicos** (Classical Errors-in-Variables) son:

Supuesto 1 (Media Cero del Error de Medición).

$$E[e] = 0 \quad (4)$$

Interpretación en este contexto:

- En promedio, los niños reportan exactamente sus horas verdaderas de TV
- Algunos sobre-reportan, otros sub-reportan, pero en promedio se cancela
- No hay sesgo sistemático hacia arriba o hacia abajo

Ejemplo numérico: Si 100 niños ven verdaderamente 10 horas/semana:

- Algunos reportarán 8, 9, 11, 12 horas
- El promedio de los reportes debe ser ≈ 10 horas
- $\frac{1}{100} \sum_{i=1}^{100} e_i \approx 0$

Supuesto 2 (Independencia del Valor Verdadero).

$$\text{Cov}(\text{tvhours}^*, e) = 0 \quad o \text{ equivalentemente} \quad E[e|\text{tvhours}^*] = 0 \quad (5)$$

Interpretación en este contexto:

- El error de medición no depende de cuánta TV realmente ve el niño
- Un niño que ve 5 horas tiene el mismo error (en expectativa) que uno que ve 20 horas
- La magnitud o dirección del error no está relacionada con el nivel verdadero de TV

Ejemplo de violación:

- Si niños que ven mucha TV (ej. 25 horas/semana) sistemáticamente subreportan más debido a vergüenza
- Entonces: $E[e|tvhours^* = 25] < 0$ y $E[e|tvhours^* = 5] \approx 0$
- Esto viola el supuesto

Supuesto 3 (Independencia de las Variables Explicativas).

$$Cov(X_j, e) = 0 \quad \text{para } j = 1, 2, \dots, k \quad (6)$$

o equivalentemente:

$$E[e|age, age^2, motheduc, fatheduc, sibs] = 0 \quad (7)$$

Interpretación específica en este contexto:

1. $\text{Cov}(\text{age}, e) = 0$:

- El error de medición no varía sistemáticamente con la edad del niño
- Niños de 8 años cometan el mismo tipo de errores que niños de 14 años

2. $\text{Cov}(\text{age}^2, e) = 0$:

- El error no tiene una relación cuadrática con la edad

3. $\text{Cov}(\text{motheduc}, e) = 0$ y $\text{Cov}(\text{fatheduc}, e) = 0$:

- El error de medición no depende de la educación de los padres
- Padres con educación universitaria reportan con la misma precisión que padres con educación secundaria
- No hay diferencias en sesgo social relacionadas con educación

4. $\text{Cov}(\text{sibs}, e) = 0$:

- El error no varía con el número de hermanos
- Familias grandes y pequeñas reportan con igual precisión

Supuesto 4 (Homocedasticidad del Error de Medición (opcional pero común)).

$$Var(e|X, tvhours^*) = \sigma_e^2 \quad (8)$$

Interpretación:

- La variabilidad del error es constante
- No depende de las características del niño ni del nivel verdadero de TV
- Todos tienen la misma "dispersión." en su error de reporte

Implicaciones de los supuestos CEV

Bajo los supuestos CEV, ¿qué ocurre con la estimación?

Cuando el error de medición está en la **variable dependiente Y** (como en este caso), los supuestos CEV implican:

Modelo verdadero:

$$\text{tvhours}^* = X\beta + u \quad (9)$$

Modelo estimado (con error de medición):

$$\text{tvhours} = \text{tvhours}^* + e = X\beta + u + e = X\beta + v \quad (10)$$

donde $v = u + e$ es el nuevo término de error.

Propiedades de v :

1. Media cero:

$$E[v] = E[u + e] = E[u] + E[e] = 0 + 0 = 0 \quad \checkmark \quad (11)$$

2. Independencia de X (supuesto clave de Gauss-Markov):

$$E[v|X] = E[u + e|X] \quad (12)$$

$$= E[u|X] + E[e|X] \quad (\text{por supuesto 3 de CEV}) \quad (13)$$

$$= 0 + E[e] \quad (\text{por supuesto 1 de CEV}) \quad (14)$$

$$= 0 + 0 = 0 \quad \checkmark \quad (15)$$

3. Varianza aumentada:

$$\text{Var}(v) = \text{Var}(u + e) \quad (16)$$

$$= \text{Var}(u) + \text{Var}(e) + 2\text{Cov}(u, e) \quad (17)$$

$$= \sigma_u^2 + \sigma_e^2 + 0 \quad (\text{si } u \text{ y } e \text{ son independientes}) \quad (18)$$

$$= \sigma_u^2 + \sigma_e^2 > \sigma_u^2 \quad (19)$$

Resultado Principal

Conclusión: Bajo los supuestos CEV, cuando el error de medición está en Y :

1. Los estimadores MCO de β son **INSESGADOS**:

$$E[\hat{\beta}] = \beta \quad (20)$$

2. Pero tienen **MAYOR VARIANZA**:

$$\text{Var}(\hat{\beta}) = \sigma_v^2(X'X)^{-1} = (\sigma_u^2 + \sigma_e^2)(X'X)^{-1} > \sigma_u^2(X'X)^{-1} \quad (21)$$

3. Los errores estándar estimados son **más grandes**
4. Los tests de hipótesis tienen **menor potencia**
5. Los intervalos de confianza son **más amplios**

En resumen: Perdemos **eficiencia** pero no **insesgamiento**.

Pregunta c: Derivación formal del sesgo por error de medición

Setup del problema

Modelo poblacional verdadero:

$$\text{tvhours}^* = \beta_0 + \beta_1 \text{age} + \beta_2 \text{age}^2 + \beta_3 \text{motheduc} + \beta_4 \text{fatheduc} + \beta_5 \text{sibs} + u \quad (22)$$

Modelo de error de medición:

$$\text{tvhours} = \text{tvhours}^* + e \quad (23)$$

Modelo que estimamos en la práctica:

$$\text{tvhours} = \beta_0 + \beta_1 \text{age} + \beta_2 \text{age}^2 + \beta_3 \text{motheduc} + \beta_4 \text{fatheduc} + \beta_5 \text{sibs} + \varepsilon \quad (24)$$

donde ε es el término de error en la regresión observada.

Paso 1: Relacionar el modelo observado con el verdadero

Sustituyendo la ecuación (7) en (8):

$$\text{tvhours} = \text{tvhours}^* + e \quad (25)$$

$$= (\beta_0 + \beta_1 \text{age} + \beta_2 \text{age}^2 + \beta_3 \text{motheduc} + \beta_4 \text{fatheduc} + \beta_5 \text{sibs} + u) + e \quad (26)$$

$$= \beta_0 + \beta_1 \text{age} + \beta_2 \text{age}^2 + \beta_3 \text{motheduc} + \beta_4 \text{fatheduc} + \beta_5 \text{sibs} + (u + e) \quad (27)$$

Comparando con la ecuación (9), vemos que:

$$\varepsilon = u + e \quad (28)$$

Conclusión del Paso 1: El término de error en la regresión observada es la suma del error del modelo verdadero (u) y el error de medición (e).

Paso 2: Analizar las propiedades del nuevo error

Para que los estimadores MCO sean insesgados, necesitamos:

$$E[\varepsilon|X] = 0 \quad (29)$$

donde X representa todas las variables explicativas: (age, age², motheduc, fatheduc, sibs).

Verificación:

$$E[\varepsilon|X] = E[u + e|X] \quad (30)$$

$$= E[u|X] + E[e|X] \quad (31)$$

$$= 0 + E[e|X] \quad (\text{por supuesto de exogeneidad del modelo verdadero}) \quad (32)$$

Ahora, por el **Supuesto 3 de CEV** ($\text{Cov}(X_j, e) = 0$ para toda X_j):

$$E[e|X] = E[e] = 0 \quad (\text{por Supuesto 1 de CEV}) \quad (33)$$

Por lo tanto:

$$E[\varepsilon|X] = 0 + 0 = 0 \quad \checkmark \quad (34)$$

Resultado Importante

Conclusión del Paso 2:

El supuesto de exogeneidad **se mantiene** bajo los supuestos CEV. Esto implica que:

$E[\hat{\beta}] = \beta \quad (\text{Los estimadores son INSESGADOS})$

(35)

No hay sesgo en los coeficientes cuando el error de medición está en la variable dependiente.

Paso 3: Analizar la varianza de los estimadores

Aunque no hay sesgo, la varianza de los estimadores sí se ve afectada.

Fórmula de la varianza del estimador MCO:

$$\text{Var}(\hat{\beta}|X) = \sigma_{\varepsilon}^2 (X'X)^{-1} \quad (36)$$

donde $\sigma_{\varepsilon}^2 = \text{Var}(\varepsilon) = \text{Var}(u + e)$.

Cálculo de σ_{ε}^2 :

Asumiendo que u y e son independientes (una extensión razonable de los supuestos CEV):

$$\sigma_\varepsilon^2 = \text{Var}(u + e) \quad (37)$$

$$= \text{Var}(u) + \text{Var}(e) + 2\text{Cov}(u, e) \quad (38)$$

$$= \sigma_u^2 + \sigma_e^2 + 0 \quad (39)$$

$$= \sigma_u^2 + \sigma_e^2 \quad (40)$$

Comparación:

- **Varianza sin error de medición:**

$$\text{Var}(\hat{\beta}_{\text{true}}) = \sigma_u^2(X'X)^{-1} \quad (41)$$

- **Varianza con error de medición:**

$$\text{Var}(\hat{\beta}_{\text{obs}}) = (\sigma_u^2 + \sigma_e^2)(X'X)^{-1} \quad (42)$$

Por lo tanto:

$$\boxed{\text{Var}(\hat{\beta}_{\text{obs}}) > \text{Var}(\hat{\beta}_{\text{true}})} \quad (43)$$

El error de medición infla la varianza de los estimadores.

Paso 4: Implicaciones para inferencia estadística

4.1. Errores estándar

Los errores estándar estimados serán:

$$\widehat{\text{SE}}(\hat{\beta}) = \sqrt{\hat{\sigma}_\varepsilon^2(X'X)^{-1}} \quad (44)$$

donde:

$$\hat{\sigma}_\varepsilon^2 = \frac{1}{n - k - 1} \sum_{i=1}^n \hat{\varepsilon}_i^2 = \frac{1}{n - k - 1} \sum_{i=1}^n (u_i + e_i)^2 \quad (45)$$

Como $\hat{\sigma}_\varepsilon^2$ estima $\sigma_u^2 + \sigma_e^2 > \sigma_u^2$:

$$\boxed{\widehat{\text{SE}}(\hat{\beta}) > \text{SE}(\hat{\beta}_{\text{true}})} \quad (46)$$

Los errores estándar son más grandes.

4.2. Estadísticos t

El estadístico t para probar $H_0 : \beta_j = 0$ es:

$$t_j = \frac{\hat{\beta}_j}{\widehat{SE}(\hat{\beta}_j)} \quad (47)$$

Dado que:

- $E[\hat{\beta}_j] = \beta_j$ (numerador correcto)
- $\widehat{SE}(\hat{\beta}_j)$ es más grande (denominador inflado)

Entonces:

$$|t_j| = \frac{|\hat{\beta}_j|}{\widehat{SE}(\hat{\beta}_j)} \quad \text{tiende a ser más pequeño} \quad (48)$$

Implicación: Mayor probabilidad de **no rechazar** H_0 (error Tipo II), es decir, menor **potencia** del test.

4.3. Intervalos de confianza

Un intervalo de confianza al 95 % es:

$$IC_{95\%}(\beta_j) = \left[\hat{\beta}_j - 1,96 \times \widehat{SE}(\hat{\beta}_j), \hat{\beta}_j + 1,96 \times \widehat{SE}(\hat{\beta}_j) \right] \quad (49)$$

Como $\widehat{SE}(\hat{\beta}_j)$ es más grande:

Los intervalos de confianza son MÁS AMPLIOS

(50)

Menor precisión en la estimación.

4.4. R^2 (Bondad de ajuste)

El R^2 se calcula como:

$$R^2 = 1 - \frac{\sum \hat{\varepsilon}_i^2}{\sum (y_i - \bar{y})^2} = 1 - \frac{SSR}{SST} \quad (51)$$

Con error de medición:

$$SSR = \sum (u_i + e_i)^2 > \sum u_i^2 \quad (52)$$

Por lo tanto:

$R_{\text{obs}}^2 < R_{\text{true}}^2$

(53)

El ajuste del modelo aparenta ser peor.

Demostración formal del sesgo: Caso general

Ahora generalizemos para entender qué pasa cuando el error está en una variable explicativa.

Modelo verdadero:

$$Y = \beta_0 + \beta_1 X^* + u \quad (54)$$

Error de medición en X :

$$X = X^* + e \quad (55)$$

Modelo estimado:

$$Y = \beta_0 + \beta_1 X + v \quad (56)$$

Derivación del sesgo

Sustituyendo (31) en (30):

$$Y = \beta_0 + \beta_1 X^* + u \quad (57)$$

$$= \beta_0 + \beta_1(X - e) + u \quad (58)$$

$$= \beta_0 + \beta_1 X - \beta_1 e + u \quad (59)$$

$$= \beta_0 + \beta_1 X + (u - \beta_1 e) \quad (60)$$

Comparando con (32): $v = u - \beta_1 e$

Verificar exogeneidad:

$$\text{Cov}(X, v) = \text{Cov}(X, u - \beta_1 e) \quad (61)$$

$$= \text{Cov}(X, u) - \beta_1 \text{Cov}(X, e) \quad (62)$$

$$= 0 - \beta_1 \text{Cov}(X^* + e, e) \quad (\text{si } u \text{ exógeno}) \quad (63)$$

$$= -\beta_1 [\text{Cov}(X^*, e) + \text{Var}(e)] \quad (64)$$

$$= -\beta_1 [0 + \sigma_e^2] \quad (\text{por Supuesto 2 de CEV}) \quad (65)$$

$$= -\beta_1 \sigma_e^2 \neq 0 \quad (66)$$

¡El supuesto de exogeneidad se viola!

Límite probabilístico del estimador

El estimador MCO converge a:

$$\text{plim}(\hat{\beta}_1) = \beta_1 + \frac{\text{Cov}(X, v)}{\text{Var}(X)} \quad (67)$$

$$= \beta_1 + \frac{-\beta_1 \sigma_e^2}{\text{Var}(X^* + e)} \quad (68)$$

$$= \beta_1 + \frac{-\beta_1 \sigma_e^2}{\text{Var}(X^*) + \sigma_e^2} \quad (\text{por independencia de } X^* \text{ y } e) \quad (69)$$

$$= \beta_1 \left(1 - \frac{\sigma_e^2}{\text{Var}(X^*) + \sigma_e^2} \right) \quad (70)$$

$$= \beta_1 \left(\frac{\text{Var}(X^*)}{\text{Var}(X^*) + \sigma_e^2} \right) \quad (71)$$

Definiendo el **ratio de confiabilidad**:

$$\lambda = \frac{\text{Var}(X^*)}{\text{Var}(X^*) + \sigma_e^2} = \frac{\text{Var}(X^*)}{\text{Var}(X)} \in (0, 1) \quad (72)$$

Entonces:

$$\boxed{\text{plim}(\hat{\beta}_1) = \lambda \beta_1} \quad (73)$$

Interpretación:

- Como $0 < \lambda < 1$, el estimador está **atenuado** (sesgado hacia cero)
- Cuanto mayor sea σ_e^2 (más error de medición), menor será λ y mayor el sesgo
- Esto se conoce como **sesgo de atenuación** (attenuation bias)

Resumen de resultados

Situación	Error en Y	Error en X
Sesgo en $\hat{\beta}$	NO	SÍ (atenuación)
Varianza de $\hat{\beta}$	Aumenta	Aumenta
Errores estándar	Más grandes	Más grandes
R^2	Disminuye	Disminuye
Potencia de tests	Disminuye	Disminuye

Cuadro 1: Efectos del error de medición según su ubicación

Aplicación al ejercicio: tvhours

En nuestro caso:

- El error de medición está en **tvhours** (variable dependiente Y)

- Las variables explicativas (age, motheduc, etc.) se asumen medidas sin error

Conclusión:

$$\boxed{\text{NO hay sesgo en los coeficientes } \hat{\beta}_j} \quad (74)$$

Pero:

1. Los estimadores tienen mayor varianza: $\text{Var}(\hat{\beta}) = (\sigma_u^2 + \sigma_e^2)(X'X)^{-1}$
2. Los errores estándar son más grandes
3. Los tests son menos potentes (mayor probabilidad de error Tipo II)
4. Los intervalos de confianza son más amplios
5. El R^2 es más bajo

Mensaje final: El error de medición en Y no destruye la validez de la inferencia causal (los coeficientes son insesgados), pero reduce la **eficiencia estadística** (precisión de las estimaciones).